

OPTIMUM UPDATE STEP FOR MOTION-COMPENSATED LIFTED WAVELET CODING

(Invited Paper)

Bernd Girod, Sangeun Han, and Chuo-Ling Chang
Information Systems Laboratory
Stanford University, Stanford, CA 94305
{bgirod, sehan, chuoling}@stanford.edu

Abstract—In motion-compensated lifted wavelet video coding, the update step typically reverses the motion vectors from the prediction step. Where motion compensation is not invertible, heuristic rules are used for the update step. This paper derives a closed-form expression for the optimal update step for a given general linear prediction step and applies the result to motion-compensated wavelet coding. Our analysis justifies using reversed motion vectors, where possible, and yields new results on proper inversion of subpixel motion compensation and optimal treatment of motion discontinuities. The analysis is extended to longer filters such as biorthogonal 5/3 wavelet transform to show the applicability of the proposed method. Experimental results confirm that optimizing the update step improves the rate-distortion performance for a motion-compensated wavelet video coder and provide justification for heuristics used in conventional techniques.

I. INTRODUCTION

Three-dimensional subband coding (3D-SBC) of video sequences using wavelet transforms can provide superior support for embedded, rate-scalable signal representations, often a requirement in best-effort networks. In an early attempt to incorporate motion compensation into 3D-SBC, Ohm [1] treats disconnected pixels arising from a non-homogeneous motion field differently from connected pixels to maintain invertibility in integer-pel accurate motion compensation. In that work, the motion vector field is restricted, and the temporal filter is limited to a two-tap Haar wavelet, which is unsatisfactory. The technique of *motion-compensated lifting* incorporates unrestricted motion compensation into 3D-SBC in a reversible fashion [2] [3]. Lifting is a procedure to implement discrete wavelet transforms [4]. Because the lifting decomposition is easily invertible, any type of operation, linear or non-linear, can be incorporated into the prediction and update steps. The lifting implementation of the wavelet transform allows for a motion-compensated temporal transform, based on any wavelet kernel and any motion model, without sacrificing the perfect reconstruction property.

Motion compensation in the prediction step should minimize the bit-rate required to encode the temporal high-band pixel. Since this bit-rate is monotonically related to

the energy of the high-band signal, this energy can be used instead to find the best motion vectors. This is analogous to minimizing residual error energy in motion-compensated predictive coding. The appropriate motion compensation for the update step, however, is not obvious. Various methods have been compared to compute backward motion fields for the update step [5], and some authors (e.g., [6]) have even reported that the update step degrades rate-distortion performance and should be omitted altogether, leading to a “truncated wavelet transform.” Typically, however, reversed motion vectors are used, in combination with heuristics for “unconnected” and “multiply connected” pixels [7] [8]. In [9], the authors provide a theoretical analysis for optimizing the update step with integer-pel accurate motion compensation. They propose update steps that average multiply connected pixels and also discuss non-linear update steps.

In [10], a closed-form expression for the optimal update step for a given general linear prediction step is derived and the result is applied to motion-compensated temporal Haar wavelet coding. The presented analysis justifies using reversed motion vectors, where possible, and yields new results on proper inversion of subpixel motion compensation and optimal treatment of motion discontinuities. In this paper the analysis is extended to longer filters such as biorthogonal 5/3 wavelet transform to show the applicability of the previously proposed method. Section II summarizes our derivation of the optimal update step for a given motion-compensated prediction step. Section III provides examples describing various cases of motion vector fields with Haar transform and Section IV extends the analysis to biorthogonal 5/3 transform. Section V presents rate-distortion comparisons for encoding video sequences with a motion-compensated wavelet coder.

II. OPTIMUM PREDICTION AND UPDATE STEPS

In [10], we consider the decomposition of a video sequence into a temporal low-band and a temporal high-band signal, using a lifting implementation as shown in Fig. 1.

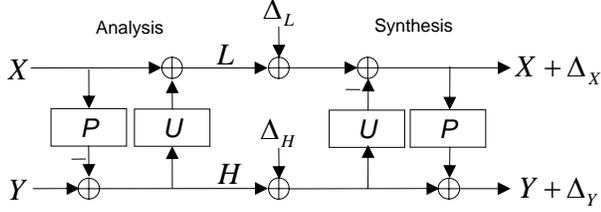


Fig. 1. Low-band/high-band decomposition by lifting.

The signals X and Y are the odd frames and the even frames of a video sequence, respectively. For our analysis, it is convenient to define X and Y as column vectors that contain all of the N pixel values in the respective odd-frame and even-frame-sequences, for example, frame by frame in line-scan order. We assume that both prediction and update steps involve exclusively linear operations, hence they can be represented by premultiplying the signal vector by a matrix. The temporal high-band signal vector H is produced by $H = Y - PX$ where P is the $N \times N$ prediction matrix. The update step $L = X + UH$ multiplies H with the $N \times N$ update matrix U to yield the low-band signal L .

We think of the prediction step P as representing motion-compensated prediction that strives to minimize the bit-rate required to encode H along with the motion vectors used for prediction. The bit-rate required to encode the low-band signal L is typically the same as would be needed to encode X directly. Since the variance of H is so much smaller than the variance of X , the exact choice of the update step U has usually only a very small impact on the bit-rate required for L . However, inspection of the inverse transform (Fig. 1) reveals that U greatly impacts the distortion in the reconstructed even and odd frame sequences $X + \Delta_X$ and $Y + \Delta_Y$. We denote by Δ_L and Δ_H the quantization errors introduced by lossy source coding, and by

$$\begin{pmatrix} \Delta_X \\ \Delta_Y \end{pmatrix} = \begin{pmatrix} I \\ P \end{pmatrix} \Delta_L + \begin{pmatrix} -U \\ I - PU \end{pmatrix} \Delta_H \quad (1)$$

the resulting errors in the reconstructed frames.

We desire to choose U such that it minimizes the mean-squared error

$$D = E\{\Delta_X^T \Delta_X + \Delta_Y^T \Delta_Y\} \quad (2)$$

We may assume that Δ_L and Δ_H are uncorrelated random vectors. Then, the choice of U does not affect the part of the distortion D due to Δ_L and we only have to consider the contribution of the high-band error, i.e.,

$$\begin{pmatrix} \Delta_X \\ \Delta_Y \end{pmatrix} = \begin{pmatrix} -U \\ I - PU \end{pmatrix} \Delta_H \quad (3)$$

Consider the autocorrelation matrix

$$\begin{aligned} R &= E \left\{ \begin{pmatrix} \Delta_X \\ \Delta_Y \end{pmatrix} \begin{pmatrix} \Delta_X \\ \Delta_Y \end{pmatrix}^T \right\} \\ &= \begin{pmatrix} -U \\ I - PU \end{pmatrix} R_H \begin{pmatrix} -U \\ I - PU \end{pmatrix}^T \end{aligned} \quad (4)$$

where $R_H = E\{\Delta_H \Delta_H^T\}$ is the autocorrelation matrix of Δ_H . We find the optimum U by matrix calculus

$$\begin{aligned} d/dU(D) &= d/dU(\text{tr}(R)) \\ &= 2(I + P^T P)UR_H - 2P^T R_H = 0 \end{aligned} \quad (5)$$

where $d/dU(\cdot)$ is a matrix whose (i, j) element is the derivative of the argument with respect to the (i, j) element of U . Assuming that R_H is full rank, we find the update matrix

$$U = (I + P^T P)^{-1} P^T \quad (6)$$

corresponding to a local extremum of the mean-squared error D (2). We note that $(I + P^T P)$ is positive definite and therefore always invertible.

Once we find the prediction step P that minimizes the bit-rate, we can easily determine the corresponding update step that minimizes the resulting distortion using (6). Except that we require the prediction step to be linear, there are no constraints on P . In the following sections, we study the case of different temporal wavelet transforms in more detail.

III. MOTION-COMPENSATED HAAR TRANSFORM

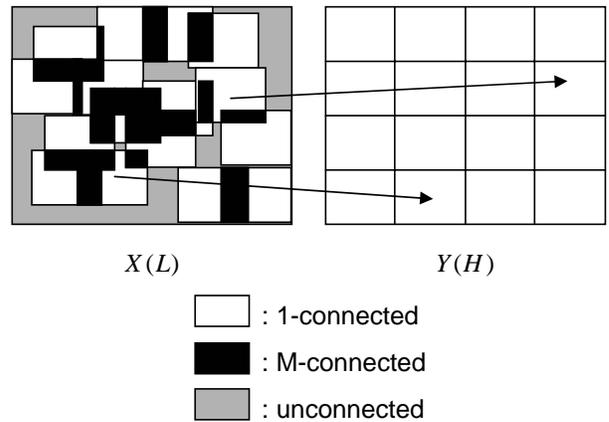


Fig. 2. Block-based motion-compensated prediction with integer displacement vectors leads to 1-connected, M-connected, and unconnected pixels in frame X .

When considering a Haar wavelet, the vectors X and Y are simply two successive frames since the Haar basis vectors do not overlap. After motion-compensated lifting, the high-band signal H corresponds to frame Y , and the

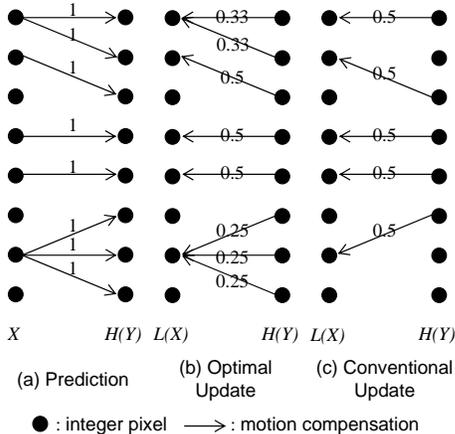


Fig. 3. Example: optimal update for integer-pel accuracy motion compensation. $L(X)$ and $H(Y)$ denote temporal references for temporal low/high-bands.

low-band signal L corresponds to X . With block-wise motion compensation, as illustrated in Fig. 2, the frame Y is divided into blocks illustrating a fixed grid, and each block finds a best match in frame X . If only integer displacements are allowed, most pixels in frame X are connected to one pixel in Y (1 -connected). However, due to the spatial variation of the displacement vector, some pixels in frame X may be used for prediction more than once ($multi$ -connected, or M -connected pixels), while others may not be used at all ($unconnected$ pixels).

Fig. 3 illustrates the optimal update step for a one-dimensional example of integer-pixel-accurate motion compensation. The circles are pixels and the arrows are the direction of motion compensation in the lifting steps. In Fig. 3, (a) is the prediction step while (b) is the optimal update step based on the (6) in Section II. The optimal solution illustrates three simple rules which are easily derived from (6):

- If a pixel in the reference frame X is 1 -connected, the corresponding pixel in the high-band H is added to the pixel along the reversed motion vectors with a weight of $\frac{1}{2}$.
- If a pixel in X is $unconnected$, this pixel is simply copied to the corresponding position in L .
- If a pixel is M -connected, all the connected pixels in the high-band H are added to the pixel with a weight of $\frac{1}{M+1}$. 1 -connected pixels are included as the special case $M = 1$.

Fig. 3(c) shows the best heuristic update known to us prior to this work [7]. We shall refer to it as the “conventional update step.” For integer-pixel accuracy, 1 -connected and $unconnected$ pixels are treated in the same way as in the optimal case. However, M -connected pixels are treated in a different way. For instance, the first encountered pixel

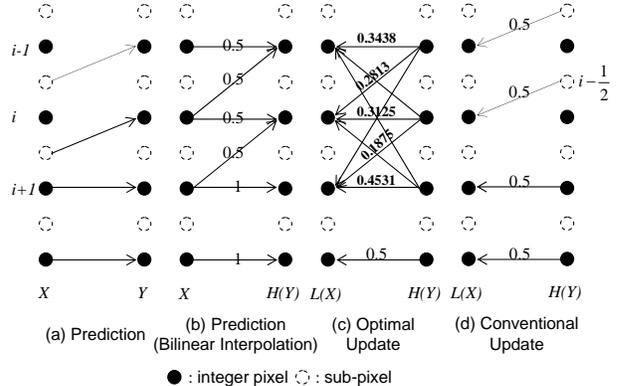


Fig. 4. Example: optimal update for half-pel accuracy motion compensation. $L(X)$ and $H(Y)$ denote temporal references for temporal low/high-bands. (In (c), only weights with absolute values > 0.1 are shown.)

in Y that uses the current M -connected pixel in X as a predictor can be chosen for computing the low-band L . Once a pixel has been selected for the update step, the same weight of $\frac{1}{2}$ is used.

A one-dimensional example of half-pixel-accurate motion compensation is shown in Fig. 4. Bilinear interpolation for sub-pixel positions is used. The resulting weights needed for the matrix P are shown in (b). Pixels in both frames X and Y might be M -connected. However, there might still be 1 -connected pixels and $unconnected$ pixels, as in the integer-pixel-accuracy case, and their optimal update step is the same as in that case. However, M -connected pixels can result in many more pixels involved in generating the optimum low-band signal at these pixel locations. In optimum update step (c), only weights with absolute values greater than 0.1 are shown. Weights not shown are not necessarily zero.

In Fig. 4(d), the conventional update step is shown [8]. When the motion displacement points to a sub-pixel position in the frame X such as $X(i + \frac{1}{2})$, conventional techniques default to the nearest integer-pixel position $X(i)$, and use $H(i - \frac{1}{2})$ for the update step. Note that bilinear interpolation in the high-band yields similar weights to those in the optimal update.

IV. MOTION-COMPENSATED BIORTHOGONAL 5/3 TRANSFORM

The entire development in Sec. II is derived for arbitrary wavelet kernel. However, in the case of longer and overlapping wavelet kernels, the proposed method should be modified in order to preserve compact support for matrices P and U . In this section, we consider the biorthogonal 5/3 wavelet transform as shown in Fig. 5. The signal X_i is each frame of a video sequence. Similar to our general case analysis, X_i is defined as a column vector that contains all

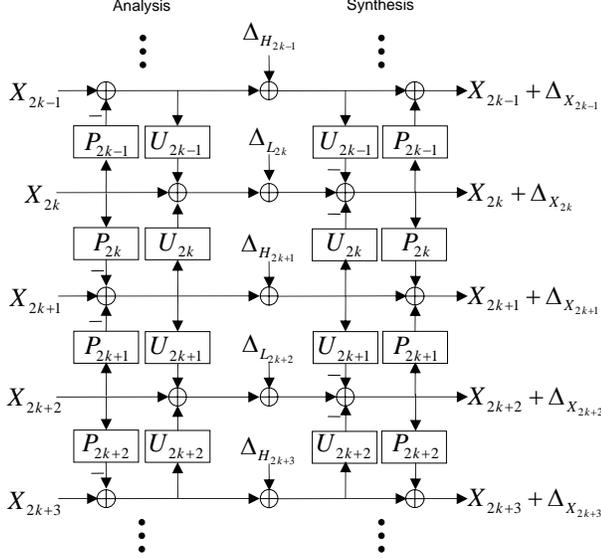


Fig. 5. Low-band/high-band decomposition by biorthogonal 5/3 lifting.

of the N pixel values in the corresponding frame. With the $N \times N$ prediction and update matrices P_i and U_i , the resulting temporal high-band and low-band signal vectors are

$$H_{2k+1} = X_{2k+1} - P_{2k}X_{2k} - P_{2k+1}X_{2k+2} \quad (7)$$

$$L_{2k} = X_{2k} + U_{2k-1}H_{2k-1} + U_{2k}H_{2k+1} \quad (8)$$

We may assume that Δ_{L_i} and Δ_{H_i} are uncorrelated, and Δ_{H_i} and Δ_{H_j} are uncorrelated. Then, the choice of U_{2k} and U_{2k+1} does not affect the part of the distortion D due to other than $\Delta_{H_{2k+1}}$ and we only have to consider the contribution of $\Delta_{H_{2k+1}}$. Without the loss of generality, quantization errors in the other temporal subband are set to zero, i.e.,

$$\Delta_X = \begin{pmatrix} \Delta_{X_{2k-1}} \\ \Delta_{X_{2k}} \\ \Delta_{X_{2k+1}} \\ \Delta_{X_{2k+2}} \\ \Delta_{X_{2k+3}} \end{pmatrix} = \begin{pmatrix} -P_{2k-1}U_{2k} \\ -U_{2k} \\ I - P_{2k}U_{2k} - P_{2k+1}U_{2k+1} \\ -U_{2k+1} \\ -P_{2k+2}U_{2k+1} \end{pmatrix} \Delta_{H_{2k+1}} \quad (9)$$

Let

$$Z = \begin{pmatrix} -P_{2k-1}U_{2k} \\ -U_{2k} \\ I - P_{2k}U_{2k} - P_{2k+1}U_{2k+1} \\ -U_{2k+1} \\ -P_{2k+2}U_{2k+1} \end{pmatrix}, \quad (10)$$

then

$$\Delta_X = Z\Delta_{H_{2k+1}}. \quad (11)$$

The mean-squared error due to $\Delta_{H_{2k+1}}$ is given by

$$D = E\{\Delta_X^T \Delta_X\} \quad (12)$$

Similarly, consider the autocorrelation matrix

$$R = E\{\Delta_X \Delta_X^T\} = Z R_{H_{2k+1}} Z^T \quad (13)$$

where $R_{H_{2k+1}}$ is the autocorrelation matrix of $\Delta_{H_{2k+1}}$. To choose U_{2k} and U_{2k+1} such that they minimize the mean-squared error D (12), we take a derivative of D with respect to U_{2k} and U_{2k+1} and set it zero as in (14) and (15).

Assuming that $R_{H_{2k+1}}$ is full rank, we find the update matrices U_{2k} and U_{2k+1} given in (16) and (17). In case of biorthogonal 5/3 wavelet transform, therefore, U can be solved with 4 neighboring P s. Note that the update weights of $\frac{1}{4}$ used typically in the 5/3 wavelet transform are not optimal. The best value of $U = \frac{2}{7}$ is found by substituting $P = \frac{1}{2}$ for all prediction matrices in (16) and (17).

V. EXPERIMENTAL RESULTS

The experiments in this section compare the optimum update step given by (6) in Section II with the conventional update step [8]. Fig. 6 shows the encoder of the wavelet video coding system with motion-compensated lifting. We use variable-blocksize motion estimation [11] and furthermore, we employ fractional-pel accuracy motion compensation with bilinear interpolation. After temporal decomposition, to further exploit the coherence among neighboring pixels within each temporal subband, a multi-level 2-D spatial DWT is then applied to decompose the subband into wavelet coefficients. Finally, the SPIHT (Set Partitioning in Hierarchical Trees) [12] algorithm is used to encode the wavelet coefficients of each temporal subband into a scalable bitstream. The SPIHT algorithm provides an embedded representation so that rate-distortion curves can be obtained by simply truncating the coded bitstreams.

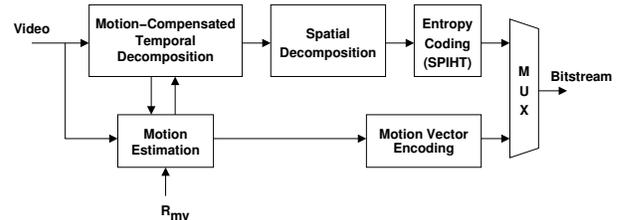


Fig. 6. Encoder structure of scalable interframe wavelet video coding system

For the optimal update scheme, the prediction matrix P can be populated from the block-wise motion vectors. Like

$$\begin{aligned} d/dU_{2k}(D) &= d/dU_{2k}(tr(R)) \\ &= ((I + P_{2k}^T P_{2k} + P_{2k-1}^T P_{2k-1})U_{2k} + P_{2k}^T P_{2k+1}U_{2k+1} - P_{2k}^T)2R_{HH_{2k+1}} = 0 \end{aligned} \quad (14)$$

$$\begin{aligned} d/dU_{2k+1}(D) &= d/dU_{2k+1}(tr(R)) \\ &= ((I + P_{2k+1}^T P_{2k+1} + P_{2k+2}^T P_{2k+2})U_{2k+1} + P_{2k+1}^T P_{2k}U_{2k} - P_{2k+1}^T)2R_{HH_{2k+1}} = 0 \end{aligned} \quad (15)$$

$$\begin{aligned} U_{2k} &= (I + P_{2k}^T P_{2k} + P_{2k-1}^T P_{2k-1} - P_{2k}^T P_{2k+1}(I + P_{2k+1}^T P_{2k+1} + P_{2k+2}^T P_{2k+2})^{-1} P_{2k+1}^T P_{2k})^{-1} \\ &\quad (P_{2k}^T - P_{2k}^T P_{2k+1}(I + P_{2k+1}^T P_{2k+1} + P_{2k+2}^T P_{2k+2})^{-1} P_{2k+1}^T) \end{aligned} \quad (16)$$

$$\begin{aligned} U_{2k+1} &= (I + P_{2k+1}^T P_{2k+1} + P_{2k+2}^T P_{2k+2} - P_{2k+1}^T P_{2k}(I + P_{2k}^T P_{2k} + P_{2k-1}^T P_{2k-1})^{-1} P_{2k}^T P_{2k+1})^{-1} \\ &\quad (P_{2k+1}^T - P_{2k+1}^T P_{2k}(I + P_{2k}^T P_{2k} + P_{2k-1}^T P_{2k-1})^{-1} P_{2k}^T) \end{aligned} \quad (17)$$

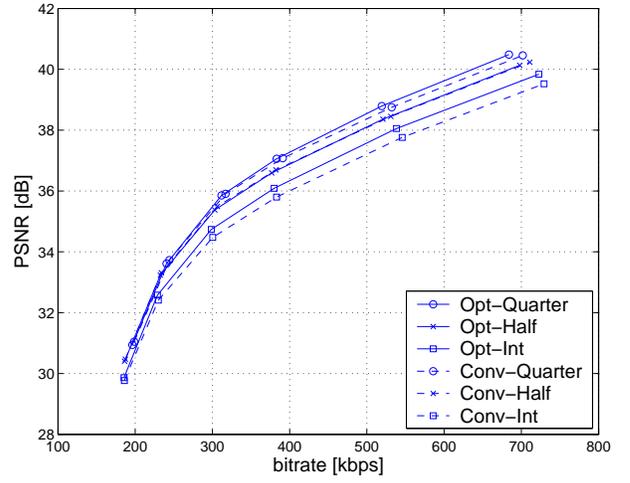
the encoder, the decoder only needs the motion vectors to construct P , so this matrix P need not be transmitted. P is a very large, but sparse matrix; e.g., 25344×25344 for QCIF sequences (176×144 pixels) in the case of temporal Haar transform. Obtaining the optimum update matrix U is challenging since the inverse of a sparse matrix is not necessarily sparse. In our implementation, we therefore do not calculate U explicitly. Instead, we observe that only UH is needed to produce the temporal low-band L . Our computation then proceeds as follows:

- Set $Z = UH = (I + P^T P)^{-1} P^T H$
- We note $(I + P^T P)Z = P^T H$.
- Solve $AZ = B$ where $A = I + P^T P$ and $B = P^T H$, exploiting the sparseness of P [13].
- Add Z to X .

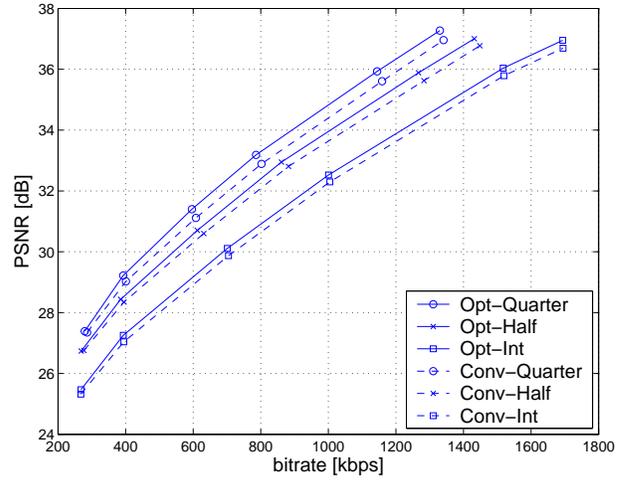
Fig. 7 shows luminance PSNR over the total bitrate for two test sequences, *Foreman* and *Mobile & Calendar*, both consisting of 288 frames in QCIF format, encoded with three levels of temporal decomposition with the motion-compensated Haar transform. We use the same motion accuracy of integer, half and quarter-pel at all decomposition levels. We observe that the optimal update step performs only slightly better than the conventional update step. This justifies the heuristics used in the conventional update step [8].

VI. CONCLUSION

We have derived a closed-form expression for the optimal update step for a given general linear prediction step and applied the result to motion-compensated wavelet coding. Our analysis provides justification for using reversed motion vectors, where possible. Unconnected and 1-connected pixels are treated as in the conventional update scheme. However, multiply-connected pixels can result in numerous pixels with different weights involved in generating the best low-band signal. We also extended the analysis to the biorthogonal $5/3$ wavelet transform and showed the applicability of the proposed method to



(a) *Foreman*



(b) *Mobile & Calendar*

Fig. 7. Rate-Distortion performance with optimal and conventional update step. Integer, half, quarter-pel accurate motion compensations are used.

longer and overlapping wavelet kernels. A new sparse matrix technique is presented that allows the practical implementation of the optimal update step for motion-compensated wavelet video coding. Experimental results show that the optimal update outperforms the conventional update step by at most 0.4dB, and that conventional update using the reversed motion vectors does nearly as well as the optimal update, thereby justifying the heuristics used in the conventional method.

REFERENCES

- [1] J.-R. Ohm, "Three-dimensional subband coding with motion compensation," *IEEE Transactions on Image Processing*, vol. 3, no. 5, pp. 559–571, Sept. 1994.
- [2] A. Secker and D. Taubman, "Motion-compensated highly scalable video compression using an adaptive 3D wavelet transform based on lifting," in *Proceedings of the IEEE International Conference on Image Processing (ICIP)*, Thessaloniki, Greece, Oct. 2001.
- [3] B. Pesquet-Popescu and V. Bottreau, "Three-dimensional lifting schemes for motion compensated video compression," in *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Salt Lake City, UT, May 2001.
- [4] W. Sweldens, "The lifting scheme: A construction of second generation wavelets," *SIAM Journal on Mathematical Analysis*, vol. 29, no. 2, pp. 511–546, 1998.
- [5] N. Božinović, J. Konrad, T. André, M. Antonini, and M. Barlaud, "Motion-compensated lifted wavelet video coding : Toward optimal motion/transform configuration," in *12th European Signal Processing Conference*, Vienna, Austria, Sept. 2004, To appear.
- [6] L. Luo, S. Li, Z. Zhuang, and Y.-Q. Zhang, "Motion compensated lifting wavelet and its application in video coding," in *Proceedings of the IEEE International Conference on Multimedia and Expo (ICME)*, Tokyo, Japan, Oct. 2001.
- [7] V. Bottreau, E. Barrau, and A. Bourge, "Architecture and features of a fully scalable motion-compensated 3d subband codec," ISO/IEC JTC1/SC29/WG11, MPEG2002/M7977, Jeju, Korea, March 2002.
- [8] J. W. Woods and P. Chen, "Improved MC-EZBC with quarter-pixel motion vectors," ISO/IEC JTC1/SC29/WG11, MPEG2002/M8366, Fairfax, VA, May 2002.
- [9] C. Tillier, B. Pesquet-Popescu, and M. van der Schaar, "Improved update operators for lifting-based motion-compensated temporal filtering," *IEEE Signal Processing Letters*, in print.
- [10] B. Girod and S. Han, "Optimum update for motion-compensated lifting," *IEEE Signal Processing Letters*, in print.
- [11] ITU-T Rec. H.264 / ISO/IEC 11496-10, "Advanced Video Coding, Final Committee Draft," Document JVT-E022, Sept. 2002.
- [12] A. Said and W. A. Pearlman, "A new fast and efficient image codec based on Set Partitioning in Hierarchical Trees," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 6, no. 3, pp. 243–250, June 1996.
- [13] C. C. Paige and M. A. Saunders, "LSQR: An algorithm for sparse linear equations and sparse least squares," *ACM Transactions on Mathematical Software*, vol. 8, no. 1, pp. 43–71, Mar. 1982.